

Statistical Machine Learning (W4400)

Fall 2015

<https://courseworks.columbia.edu>

John P. Cunningham

jpc2181@columbia.edu

Xiaoou Li

xiaoou@stat.columbia.edu

Chris Mulligan

clm2186@columbia.edu

Homework 4

Due: Tuesday 24 November 2015

Students are encouraged to work together, but homework write-ups must be done individually and must be entirely the author's own work. Homework is due at the **beginning** of the class for which it is due. **Late homework will not be accepted under any circumstances.** To receive full credit, students must thoroughly explain how they arrived at their solutions and include the following information on their homeworks: name, UNI, homework number (e.g., HW03), and class (STAT W4400). All homework must be turned in online through Courseworks in PDF format, have a .pdf extension (not zip or other archive!), and be less than 4MB. If programming is part of the assignment, the code must be turned in in one or more .R files. Homeworks not adhering to these requirements will receive no credit.

1. PCA (15 points)

An application of PCA popular in Computer Vision are "eigenfaces", where PCA is applied to a face image data base, typically for use in face recognition algorithms. This method is very similar to the digit example we have seen in class (which could analogously be named "eigendigits", and somebody somewhere probably has done so).

The figure below shows (on the left) 100 examples from a data set of face images. Each image is a 92×112 grayscale image and can be interpreted as a vector $x \in \mathbb{R}^{10304}$. The figure on the right shows the first 48 principal components ξ_1, \dots, ξ_{48} , visualized as images.



1. (5 points) How many principal components are there in total?
2. (10 points) Describe a method that *approximately* reconstructs a specific face image $x \in \mathbb{R}^{10304}$ using only the first 48 principal components? To describe your method, denote by \hat{x} the approximate representation of x . Represent \hat{x} as an equation $\hat{x} = \dots$. Please define precisely what each variable occurring on the right-hand side of the equation means.

2. k -means (35 points)

In this exercise, you will prove that the k -means algorithm minimizes a sensible clustering objective. To begin, recall the k -means setup: we are given a dataset $x_1, \dots, x_n \in \mathbf{R}^d$ and a desired number of clusters K , and we seek a cluster assignment $m_i \in \{1, \dots, K\}$ of each data point x_i .

Homework problems.

- (5 points) Define the set $C_k = \{i | m_i = k\}$. What does C_k describe in words? Let $|C_k|$ be the number of elements in the set C_k .
- (5 points) Consider the objective:

$$f(m_1, \dots, m_n) = \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j \in C_k} \|x_i - x_j\|_2^2.$$

Describe this objective in words. Why is this a sensible clustering objective?

- (15 points) Define $\mu_k = \frac{1}{|C_k|} \sum_{j \in C_k} x_j$. Prove the following:

$$\frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j \in C_k} \|x_i - x_j\|_2^2 = 2 \sum_{i \in C_k} \|x_i - \mu_k\|_2^2.$$

- (10 points) Prove that the previous result implies that the k -means algorithm decreases the objective $f(m_1, \dots, m_n)$ with each step of the algorithm.

3. Multinomial Clustering (50 points)

In this exercise, we will apply the EM algorithm and a finite mixture of multinomial distributions to the image segmentation problem. Image segmentation refers to the task of dividing an input image into regions of pixels that “belong together” (which is about as formal and precise a definition as you will be able to find in the literature.) One way to approach image segmentation is to extract a suitable set of features from the image and apply a clustering algorithm to them. The clusters output by the algorithm can then be regarded as the image segments. The features we will use are histograms drawn from small regions in the image.

Histogram clustering: The term *histogram clustering* refers to grouping methods whose input features are histograms. A histogram can be regarded as a vector; if all input histograms have the same number of bins d , a set of input histograms can be regarded as a set of vectors in \mathbb{R}^d . We can therefore perform a simple histogram clustering by running a k -means algorithm on this set of vectors. The method we will implement in this problem is slightly more sophisticated: Since histograms are represented by multinomial distributions, we describe the clustering problem by a finite mixture of multinomial distributions, which we estimate with the EM algorithm.

The EM algorithm for multinomial mixtures. The input of the algorithm is:

- A matrix of histograms, denoted \mathbf{H} , which contains one histogram vector in each row. Each column corresponds to a histogram bin.
- An integer K which specifies the number of clusters.
- A threshold parameter τ .

The variables we need are:

- The input histograms. We denote the histogram vectors by $\mathbf{H}_1, \dots, \mathbf{H}_n$.
- The *centroids* $\mathbf{t}_1, \dots, \mathbf{t}_K$. These are \mathbb{R}^d vectors (just like the input features). Each centroid is the parameter vector of a multinomial distribution and can be regarded as the center of a cluster; they are computed as the weighted averages of all features assigned to the cluster.
- The *assignment probabilities* $\mathbf{a}_1, \dots, \mathbf{a}_n$. Each of the vectors \mathbf{a}_i is of length K , i. e. contains one entry for each cluster k . The entry a_{ik} specifies the probability of feature i to be assigned to cluster k . The matrix which contains the vectors \mathbf{a}_i in its rows will be denoted \mathbf{P} .

The algorithm iterates between the computation of the assignment probabilities \mathbf{a}_i and adjusting the cluster centroids \mathbf{t}_k .

The algorithm:

1. Choose K of the histograms at random and normalize each. These are our initial centroids $\mathbf{t}_1, \dots, \mathbf{t}_K$.
2. Iterate:

(a) E-step: Compute the components of each \mathbf{a}_i , the assignment probabilities a_{ik} , as follows:

$$\begin{aligned}\phi_{ik} &:= \exp \left(\sum_j H_{ij} \log(\mathbf{t}_{kj}) \right) \\ a_{ik} &:= \frac{c_k \phi_{ik}}{\sum_{l=1}^K c_l \phi_{il}}\end{aligned}$$

Note that ϕ_{ik} is the multinomial probability of \mathbf{H}_i with parameters \mathbf{t}_k , up to the multinomial coefficient $\frac{n!}{H_{i1}! \dots H_{id}!}$, which cancels out in the computations of a_{ik} .

(b) M-step: Compute new mixture weights c_k and class centroids \mathbf{t}_k as

$$\begin{aligned} c_k &= \frac{\sum_{i=1}^n a_{ik}}{n} \\ \mathbf{b}_k &= \sum_{i=1}^n a_{ik} \mathbf{H}_i \\ \mathbf{t}_k &= \frac{\mathbf{b}_k}{\sum_{j=1}^d b_{kj}} \end{aligned}$$

(c) Compute a measure of the change of assignments during the current iteration:

$$\delta := \|\mathbf{A} - \mathbf{A}^{\text{old}}\| ,$$

where \mathbf{A} is the assignment matrix (with entries a_{ik}), \mathbf{A}^{old} denotes assignment matrix in the previous step, and $\|\cdot\|$ is the matrix 1-norm (the largest sum of column absolute values). This is implemented in R as `norm(., "0")`.

3. Terminate the iteration when $\delta < \tau$.
4. Turn the soft assignments into a vector \mathbf{m} of hard assignments by computing, for $i = 1, \dots, n$,

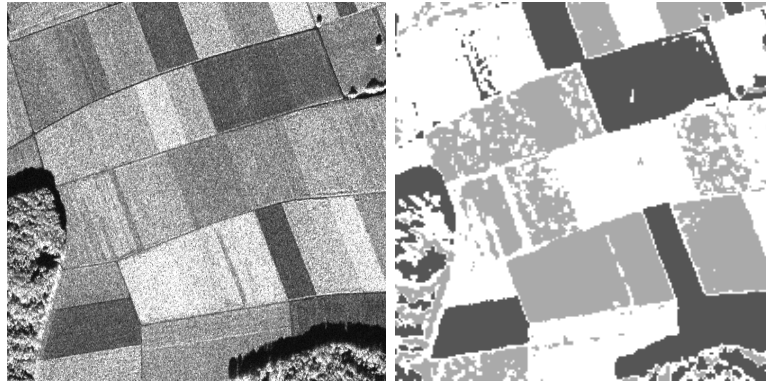
$$m_i := \arg \max_{k=1, \dots, K} a_{ik} ,$$

i. e. the index k of the cluster for which the assignment probability a_{ik} is maximal.

The histogram data: The histograms we use have been extracted from an image by the following procedure:

1. Select a subset of pixels, called the *sites*, at which we will draw a histogram. Usually, this is done by selecting the pixels at the nodes of an equidistant grid. (For example, a 2-by-2 grid means that we select every second pixel in every second row as a site.)
2. Place a rectangle of fixed radius around the site pixel.
3. Select all pixels within the rectangle and sort their intensity values into a histogram.

The data we provide you with was drawn from the 800x800 grayscale image shown below (left):



The image is a radar image; more specifically, it was taken using a technique called synthetic aperture radar (or SAR). This type of image is well-suited for segmentation by histogram clustering, because the local intensity distributions provide distinctive information about the segments. The image on the right is an example segmentation using $K = 3$ clusters, computed with the algorithm described above.

The histograms were drawn at the nodes of a 4-by-4 pixel grid. Since the image has 800×800 pixels, there are $200 \times 200 = 40000$ histograms. Each was drawn within a rectangle of edge length 11 pixels, so each histogram contains $11 \times 11 = 121$ values.

Homework problems. Download the data file `histograms.bin` from the course homepage. You can load it in R with the command

```
H<-matrix(readBin("histograms.bin", "double", 640000), 40000, 16)
```

which results in a 40000×16 matrix. Each row is a histogram with 16 bins. You can check that the matrix has been loaded correctly by saying `dim(H)` (which should be 40000 rows by 16 columns) and `rowSums(H)` (the sum of each row should be 121).

1. Implement the EM algorithm in R. The function call should be of the form `m<-MultinomialEM(H,K,tau)`, with `H` the matrix of input histograms, `K` the number of clusters, and `tau` the threshold parameter τ .
2. Run the algorithm on the input data for $K=3$, $K=4$ and $K=5$. You may have to try different values of τ to obtain a reasonable result.
3. Visualize the results as an image. You can turn a hard assignment vector `m` returned by the EM algorithm into an image using the `image` function in R.

Hints:

- You may encounter numerical problems due to empty histogram bins. If so, add a small constant (such as 0.01) to the input histograms.
- Both the E-step and the M-step involve sums, which are straightforward to implement as loops, but loops are slow in interpreted programming languages (such as R or Matlab). If you want to obtain a fast algorithm, try to avoid loops by implementing the sums by matrix/vector operations.
- Be sure to submit your visualized image results along with your code.